

Title: Measurement Science for Complex Information Systems

Submitting OU: Information Technology Laboratory

Preferred Category: competence

Champion(s): K. Mills, C. Dabrowski, V. Marbukh, F. Hunt, and J. Filliben

What is the problem and why is it hard?

Global behavior in complex information systems, such as the Internet, computing grids, and service-oriented architectures, is difficult to measure, understand and control. Subtle interactions and indirect couplings through shared resources lead to emergent behaviors, which are not predictable from analyzing the designs of individual components. Further, systems exhibit emergent phenomena over large scale, which presents significant challenges for current measurement, analysis and control regimes. As advances in information technology increasingly drive the nation's economy and bolster national defense, ensuring sustained technical innovation will depend upon the availability of a consistent and reliable metrology. This required metrology must encompass a science of complex information systems, i.e., systems where interactions exhibit coherent properties that extend over large scale. NIST is uniquely poised to make seminal contributions to a measurement science for complex information systems.

How is it solved today, and by whom?

Emergent behavior¹⁻² in information systems is largely unaddressed today because industry lacks ability to measure at large scale with sufficient timeliness. Yet, as global information systems scale, emergent behavior poses serious threats to system stability, and consequently, to economic productivity. For example, coincident oscillation of network congestion – arising from Internet congestion-control algorithms³ – can lead to sluggish response times and wave-like traffic surges. More generally, the software industry is moving to a distributed, component-based paradigm for development and deployment; thus, software components can be tested in a laboratory and unleashed to thousands of users throughout a network – only then to discover unforeseen effects from emergent behaviors. For example, updating grid clients to combat denial-of-service attacks can cause a system-wide decay in job completion times due to an emergent property of distributed, resource-allocation algorithms⁴. One can extrapolate to future cases where updated components may be distributed widely only to stimulate unexpected emergent behaviors that did not appear during design, development, and testing. Such behaviors cannot be detected currently for two main reasons: (1) designers and developers have insufficient technical approaches and supporting tools to measure and analyze global system behavior and (2) testing laboratories possess insufficient techniques to exercise systems at scales required to reveal most emergent behaviors.

The current state of the art can be discussed along two lines: industry and academe. Industry is investigating technology for autonomic computing⁵ that will allow system elements to measure the environment and adapt configuration and behavior dynamically. Autonomic computing research aims to reduce the costs of managing deployed systems. Technologies of this nature are likely to create a step increase in emergent behaviors because components must be tested for proper operation alone and in combination with other components as well as changing environmental conditions. Testing components in multiple combinations is currently a challenge⁶ – adding consideration of the environment expands the potential state space further beyond industry's ability to cope.

Academe is investigating the use of large-scale, reconfigurable test and emulation systems (such as Emulab⁷ and DETER⁸). These facilities establish controlled environments of large scale to test performance, robustness, and security in system designs. These facilities could also be used as platforms for large simulations. Expected availability of multi-core, multi-chip microprocessors⁹ presage an order of magnitude increase in the scale of such test facilities, which

would permit designers to test systems of larger size. Yet, without significant advancement in techniques to measure and analyze behavior in distributed systems, larger test facilities will be insufficient. Some academic work is underway to develop techniques to visualize and characterize network traffic. Visualization techniques aim to reveal dynamic changes in network structure as well as network behavior¹⁰⁻¹¹. Most attempts to characterize network traffic¹² consider the application of power-spectral density, $1/f$ noise, and wavelets to reveal timescales of coherent traffic patterns. Much of this work aims to understand patterns of Internet traffic and to monitor changes over time. None of these approaches can directly identify causes associated with behavioral changes nor suggest how particular behaviors can be controlled. A few researchers¹³⁻¹⁵ are investigating the application of utility functions as a means to assert decentralized control over network behavior. Overall, recent research on large, complex systems provides a critical mass of work on which to base a measurement science; however, establishing a coherent science will require additional effort.

What is the new idea? Why can we succeed now?

To form the basis for a coherent science, we propose the application of modeling and analysis techniques from the physical sciences in order to define a systematic method to measure, understand and control emergent behaviors in complex information systems. To do this, we can draw on previous work. For example, using cellular automata¹⁶⁻¹⁹ (CA), a homogeneous modeling technique, physicists and chemists represent and analyze large collections of multi-dimensional elements, such as molecules. Others apply CA to understand the behavior of landslides²⁰, earthquakes²¹, soil contamination²², and transportation systems²³. The abstract nature of CA allows them to represent systems of much larger scale than possible using more detailed models. The basic idea underlying such approaches is to model interaction phenomena at large scale; measure and analyze data generated by the model; and gain insight into the dynamics underlying the phenomena. Increasing the scale and abstraction of a model permits global behaviors to be measured more clearly. For example, behavior of abstract models may be represented as a multi-dimensional space that can be subjected to various spatial-temporal analyses²⁴⁻²⁵ to seek patterns, and to reveal causality. Capturing state spaces from more typical models of information systems proves intractable because the level of abstraction is quite low. Even if large state spaces could be captured, analyzing the information would prove quite challenging due to the volume of noise injected from myriad irrelevant details.

Existing research on complex systems in the physical sciences provides a foundation on which we could build a coherent measurement science for complex information systems. Further, the expansion of cost-effective computation, expected over the next few years, provides new opportunity to create tractable models in the large. We should now be able to construct, exercise, and analyze abstract models (we plan to focus on agent automata²⁶⁻²⁸ – a recent evolution of CA) to measure, understand, and control global behaviors in large-scale systems, such as the Internet, computing grids, and service-oriented architectures. Such models should provide a foundation for systematic analysis, and, ultimately, a measurement science for complex information systems.

Technical plan:

Development of the necessary analysis method faces three main challenges: (1) construct agent automata (AA) models of large-scale systems that yield sufficiently accurate representations of the phenomena modeled, (2) identify tractable techniques to measure system evolution and reveal causality, and (3) characterize decentralized algorithms to elicit desired behaviors and thus verify the causal analysis. To address the model-construction challenge, we envision designing, implementing, and exercising AA models from which we collect measurements that may be compared against similar measurements from more detailed models of the same phenomena and, where appropriate, from measurements collected from comparable operational phenomena. In the Internet case, detailed simulators²⁹⁻³³ exist against which we can compare measurements obtained from an Internet AA model, which we will develop. Further, there exists a growing literature on

Internet traffic measurements³⁴⁻³⁸ against which we can compare results from our Internet AA model. In the computing-grid case, there exists a growing population of detailed simulators³⁹⁻⁴¹ against which we can compare measurements from a Grid AA model, which we will develop. Presently, few measurements exist from operational grids; however, we plan to conduct comparisons between our Grid AA model and any accepted measurements that become available during the project. Validation of our AA models is a critical step because computer scientists do not currently accept that highly abstract models of large-scale systems capture behavior with sufficient fidelity to yield useful measurements. A major innovation from this project is production of validated AA models for selected large-scale (10^6 element) distributed systems.

To address the tractable-analysis challenge, we envision a human-guided approach to exploratory data analysis⁴²⁻⁴³. We will record AA evolution as a multi-dimensional space against which projections, transformations, and visualizations may be applied. Some preliminary work⁴ shows that searching for causality requires a human creative element that cannot be readily encoded in algorithms; thus, we will develop tools and techniques to support user-guided data analysis. We seek tools for comparative visualization of spatial-temporal data, as well as mathematical techniques to characterize and correlate behavioral data. As a means to measure emergent behaviors, we plan to investigate the comparative merits of known mathematical approaches to measure physical systems (e.g., self-similarity analyses⁴⁴), economic systems (e.g., correlation analyses⁴⁵), and information systems (e.g., Kolmogorov Complexity⁴⁶). As we identify effective techniques to measure behavior and to detect causality, then we can encode them as algorithms. A main innovation of our approach is to support human-guided analysis first, while automating approaches only after establishing their effectiveness. Solving the challenges of model construction and tractable analysis will allow us to define an initial method for systematic analysis of complex information systems.

To address the decentralized-control challenge, we envision investigating two main approaches: (1) market-based techniques and (2) biologically inspired mechanisms. Some academic researchers investigate the use of market-based techniques, such as auctions⁴⁷, commodity markets⁴⁸, and present-value analysis⁴⁹, to influence behavior in distributed systems. While details differ, these approaches all exploit price feedback to influence resource allocation. Other researchers investigate biologically inspired mechanisms⁵⁰⁻⁵⁴, such as morphogen gradients, chemotaxis, quorum sensing, selective stabilization, and local and lateral inhibition, to self-organize structure and behavior in distributed systems. Using our AA models and data-analysis tools, we will measure properties of various control regimes. Here, our research has two main innovations: control techniques will be (1) measured at large scale and (2) compared across disciplines (i.e., economics and biology). Solving the control challenge will verify our method and techniques for causal analysis.

Below we list key milestones for the proposed research.

- Year 1: Preliminary AA model of an Internet
- Year 1: Initial implementation of a data repository and analysis framework
- Year 2: AA model for Internet validated against detailed simulations
- Year 2: Preliminary AA model of a computing grid
- Year 2: Initial package of automated data analysis functions
- Year 3: AA model for Internet validated against measurement data
- Year 3: AA model of computing grid validated against detailed simulations
- Year 3: Report outlining initial analytic method
- Year 4: AA models augmented with economic control regimes
- Year 4: Report characterizing properties of economic control regimes
- Year 5: AA models augmented with biologically inspired control regimes
- Year 5: Report characterizing properties of biologically inspired control regimes

- Year 5: Report describing analytic method, verified by control regimes

What is the impact if successful, and who would care?

Society is increasingly reliant on complex information systems. The Internet, a key factor in commerce, education, and government, can now be coupled with component-based software. The Defense Department recently announced⁵⁵ a \$13-billion plan to reengineer its information systems as a service-oriented architecture. Similar initiatives, underway throughout government⁵⁶⁻⁵⁷ and industry⁵⁸⁻⁶⁰, will reshape the nation's information infrastructure. Creating a measurement science, and supporting technologies, to understand and control behavior in complex information systems should provide designers with the metrology needed to deploy reliable and predictable systems that meet targets for performance, robustness, and reliability; thus, customers receive better service and the economy is spared potentially large costs associated with service disruptions that lead to subsequent, repeated redesign and redeployment.

Why should NIST do this?

No science exists to measure and control tomorrow's complex information systems. The Federal government has an inherent responsibility to ensure the quality of the nation's information infrastructure, a common resource crossing boundaries of individual organizations and constituencies. Charged with improving the nation's measurement infrastructure, NIST can establish a measurement science that will enable designers to predict and control the behavior of large, complex information systems.

List Key Selling Points.

- *OSTP-OMB Priority: High-End Computing and Networking R&D* – "...research...for secure reliable, distributed computing environments...."
- *2002 NIST report: Economic Impacts of Inadequate Infrastructure for Software Testing* – "complexity of the...software needed to support the U.S.'s computerized economy is increasing at an alarming rate.... Estimates of the economic costs of faulty software in the U.S. range in the tens of billions of dollars per year."
- *2003 NSF report: Grand Challenges in Information Systems* – "...complexity of large-scale information systems...approaches that of...the economy, where technicians may understand components, but can neither predict nor control the whole system." (See also reference 53)
- *2005 GCN article: Defense issues RFP for \$13b Encore II contract* – "The Net-Centric Enterprise Services program will use Web services...." Preliminary discussions with DISA and DARPA suggest potential support for NIST research to investigate unanticipated emergent behaviors in service-oriented architectures.

References

1. Bak P. (1996) How Nature Works: the science of self-organizing criticality, Copernicus, New York.
2. Holland J. H. (1998) Emergence: From Chaos to Order, Perseus Books, Cambridge, Mass.
3. Yuan J. and Mills K. (2005) “Macroscopic Dynamics in Large-Scale Data Networks”, chapter 8 in Complex Dynamics in Communication Networks, edited by Ljupco Kocarev and Gábor Vattay, Springer, ISBN 3-540-24305-4, pp. 191-212.
4. Mills K. and Dabrowski C. (2006) “Investigating Emergent Behavior in Computing Grids”, submitted to the International Conference on Dependable Systems and Networks.
5. Ganek A. and Corbi T. (2003) “The dawn of the autonomic computing era”, in the *IBM Systems Journal*, Vol. 42, No. 1, pp. 5-18.
6. Legeard B., Peureux F., Utting M. (2004) “Controlling test case explosion in test generation from B formal models”, *Software Testing, Verification and Reliability*, John Wiley, Volume 14, Issue 2, pp. 81-103.
7. White B., et al. (2002) “An integrated experimental environment for distributed systems and networks”, *SIGOPS Operating Systems Review*, ACM, Vol. 36, special issue on Cluster resource management, pp. 255-270.
8. Zhang K., et al. (2005) “Performing BGP Experiments on a Semi-realistic Internet Testbed Environment”, *Proceedings of the Second International Workshop on Security in Distributed Computing Systems*, IEEE, Vol. 2, pp. 130-136.
9. Kumar R., Zyuban V., Tullsen D. (2005) “Interconnections in Multi-Core Architectures: Understanding Mechanisms, Overheads and Scaling”, *SIGARCH Computer Architecture News*, ACM, Vol. 33, No. 2, pp. 408-419.
10. Marron J., Hernández-Campos F., Smith F. (2002) “Mice and Elephants Visualization of Internet Traffic”, in *Proceedings of 15th Conference on Computational Statistics*.
11. Lad. M, Zhang L., Massey D. (2004) “Link-Rank: a graphical tool for capturing BGP routing dynamics”, in *Proceedings of Network Operations and Management Symposium*, IEEE, Vol. 1, pp. 627-640.
12. Karagiannis T., Molle M., Faloutsos M. (2004) “Long-Range Dependence: Ten Years of Internet Traffic Modeling”, *IEEE Internet Computing*, Sept-Oct, pp. 57-64.
13. Shenker S. (1995) “Fundamental Design Issues for the Future Internet”, *Journal On Selected Areas In Communications*, IEEE, Vol. 13, No. 7, pp. 1176-1188.
14. Wolpert D., Turner K., Frank J. (1998) “Using collective intelligence to route Internet traffic”, *Proceedings of the 1998 conference on Advances in neural information processing systems*, pp. 952-958.
15. Kunniyur S. and Srikant R. (2003) “End-to-end congestion control schemes: utility functions, random losses and ecn marks”, *IEEE/ACM Transactions on Networking*, Vol. 11, No. 5, pp. 689-703.
16. Frisch U., Hasslacher B., Pomeau Y. (1986) “Lattice-Gas Automata for the Navier-Stokes Equation”, *Phys. Rev. Lett.* 56, Issue 14-7, pp. 1505–1508.
17. Margolus N., Toffoli T., Vichniac G. (1989) “Cellular-Automata Supercomputers for Fluid-Dynamics Modeling”, *Phys. Rev. Lett.* 56, Issue 16-21, pp. 1694–1696.
18. Smith M. (1994) Cellular Automata Methods in Mathematical Physics, PhD dissertation, Massachusetts Institute Of Technology, 243 pages.
19. Wolfram S. (1994) Cellular Automata and Complexity, Westview Press, 596 pages.
20. Di Gregorio S., et al. (1994) “Landslide Simulation by Cellular Automata in a Parallel Environment”, *Proceedings of the 2nd International Workshop on Massive Parallelism: Hardware, Software, and Applications*, World Scientific, pp. 392-407.
21. Wang Z. (1995) A Parallel Implementation of a Cellular Automata Based Earthquake Model, master’s thesis, Carleton University, Carleton, Canada.

22. Spezzano G. and Talia D. (1998) "Design Parallel Models of Soil Contamination by the CARPET Language," *Future Generation Computer Systems*, Vol. 13, No. 4-5, pp. 291-302.
23. Nagel K. and Schleicher A. (1994) "Microscopic Traffic Modeling on Parallel High-Performance Computers", *Parallel Computing*, Vol. 20, January pp. 125-146.
24. Haining R. (1990) Spatial Data Analysis in the Social and Environmental Sciences, Cambridge University Press, 432 pages.
25. Shumway H. and Stoffer D. (2005) Time Series Analysis and Its Applications, Springer, 582 pages.
26. Badouel E., et al. (1999) "Modeling Dynamic Agent System with Cooperating Automata", *Proc. of Parallel and Distributed Processing Techniques and Applications*, CS-REA Press, pp. 11-17.
27. Lomazova I. (2003) "Interacting Automata for Modelling Distributed Systems", *Lecture Notes in Computer Science*, Springer-Verlag, Vol. 2658, pp. 851-860.
28. Marc F., Degirmenciyan-Cartault I., Fallah-Seghrouchni, A. (2003) "Multi-Agent Planning as a Coordination Model for Self-Organized Systems", *Proceedings of the IEEE/WIC International Conference on Intelligent Agent Technology*, pp. 218-225.
29. Barodia R., et al. (1998) "Parsec: a parallel simulation environment for complex systems", *Computer*, IEEE, Oct, pp. 77-85.
30. Cowie J., Nicol D., Ogielski A. (1999) "Modeling The Global Internet", *Computing in Science and Engineering*, IEEE, Jan-Feb, pp. 42-50.
31. Fall K. (1999) "Network Emulation in the Vint/NS Simulator", *Proceedings of the Fourth IEEE Symposium on Computers and Communications*, July, p. 244-250.
32. Breslau L., et al. (2000) "Advances in Network Simulation", *Computer*, IEEE, May, pp. 59-67.
33. Floyd S. and Paxson V. (2001) "Difficulties in Simulating the Internet", *IEEE/ACM Transactions on Networking*, Vol. 9, No. 4, pp. 392-403.
34. Crovella M. and Bestavros A. (1996) "Self-similarity in World Wide Web traffic: evidence and possible causes", *Proceedings of the 1996 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, May, p.160-169.
35. Thompson K., Miller G., Wilder R. (1997) "Wide-Area Internet Traffic Patterns and Characteristics", *Network*, IEEE, Nov-Dec, pp. 10-23.
36. Klaffy K., Miller G., Thompson K. (1998) "The Nature of the Beast: Recent Traffic Measurements from an Internet Backbone", the *Proceedings of INET'98*, July.
37. Feldman A., et al. (1999) "Dynamics of IP traffic: a study of the role of variability and the impact of control", *Proceedings of the conference on Applications, technologies, architectures, and protocols for computer communication*, ACM, pp. 301-310.
38. Fraleigh C., et al. (2003) "Packet-Level Traffic Measurements from the Sprint IP Backbone", *Network*, IEEE, Nov-Dec, pp. 6-16.
39. Buyya, R. and Murshed, M. (2002) "GridSim: a toolkit for the modeling and simulation of distributed resource management and scheduling for Grid Computing," *Concurrency and Computation: Practice and Experience*, Vol. 14, pp. 1175-1220.
40. Legrand, A., Marchal, L., Casanova, H. (2003) "Scheduling Distributed Applications: The SimGrid Simulation Framework," *Proceedings of the third IEEE International Symposium on Cluster Computing and the Grid (CCGrid'03)*, Tokyo, Japan., pp. 138-145.
41. Liu, X., Xia, H., Chien, A., (2004) "Validating and Scaling the MicroGrid: A Scientific Instrument for Grid Dynamics", *Journal of Grid Computing*, Vol. 2, No. 2, pp. 141-161.
42. Han J. and Kamber M. (2000) Data Mining: Concepts and Techniques, Morgan Kaufmann Publishers, 550 pages.
43. Gelman A. (2004) "Exploratory Data Analysis for Complex Models", *Journal of Computational and Graphical Statistics*, Vol. 13, No. 4, pp. 755-779.
44. Karagiannis, T., Faloutsos M., Molle, M. (2003) "A User-Friendly Self-Similarity Analysis Tool", *Computer Communication Review*, ACM SIGCOMM, 13 pages.

45. Yuan J. and Mills K. (2005) "Cross-Correlation Based Method for Spatial-Temporal Traffic Analysis ", *Performance Evaluation*, Volume 61/2-3, pp 163-180.
46. Gammerman A. and Vovk V. (1999) "Kolmogorov Complexity: Sources, Theory and Applications" in *The Computer Journal*, Vol. 42, No. 4, pp. 252-255.
47. Waldspurger C., et al. (1992) "Spawn: A Distributed Computational Economy", *IEEE Transactions on Software Engineering*, Vol. 18, No. 2, pp. 103-117.
48. Wolski R., et al. (2003) "Grid Resource Allocation and Control Using Computational Economies", in: F. Berman, G. Fox, and A. Hey (eds.), Grid Computing: Making the Global Infrastructure a Reality, Wiley and Sons, pp. 747-772.
49. Irwin D., Grit L., Chase J. (2004) "Balancing Risk and Reward in a Market-based Task Service", *Proceedings of the 13th International Symposium on High Performance Distributed Computing*, pp. 160-169.
50. Tarasewich P. and McMullen P. (2002) "Swarm intelligence: power in numbers", *Communications of the ACM*, Vol. 45 No. 8, pp. 62-67.
51. Nagpal R. (2004) "A Catalog of Biologically-inspired Primitives for Engineering Self-Organization", in Engineering Self-Organising Systems: Nature-Inspired Approaches to Software Engineering, Springer, Vol. 2977, pp. 53-62.
52. Lodding K. (2004) "Hitchhiker's Guide to Biomorph Software", *QUEUE*, June, pp. 66-76.
53. Suzuki J. and Suda T. (2005) "A Middleware Platform for a Biologically-inspired Network Architecture Supporting Autonomous and Adaptive Applications", *IEEE Journal on Selected Areas in Communications (JSAC)*, Special Issue on Intelligent Services and Applications in Next Generation Networks, vol. 23, no. 2, pp. 249-260.
54. Babaoglu O., et al. (editors) (2005) Self-star Properties in Complex Information Systems: Conceptual and Practical Foundations, Springer, Volume 3460, 447 pages.
55. Onley, D. (2005) "DOD will spend \$13b to convert legacy systems", *Government Computer News*, Vol. 24 No. 34.
56. Executive Office of the President (2004) Expanding E-Government: Partnering for a Results-Oriented Government, Office of Management and Budget, December.
57. Harris B. (2006) "Trends in Enterprise Architecture", *Government Technology*, January 10.
58. Gardner G., Baumstarck J., Sego P. (2005) Industry Best Practices in Achieving Service Oriented Architecture (SOA), Association for Enterprise Integration, April 22, 31 pages.
59. Bieberstein N., et al. (2005) Service-Oriented Architecture (SOA) Compass: Business Value, Planning, and Enterprise Roadmap, IBM Press, 272 pages.
60. Tabor Communications (2006) The Emergence of Grid and Service-Oriented IT: An Industry Vision for Business Success, to appear in February.